



Dynamic R-CNN: Towards High Quality Object Detection via Dynamic Training

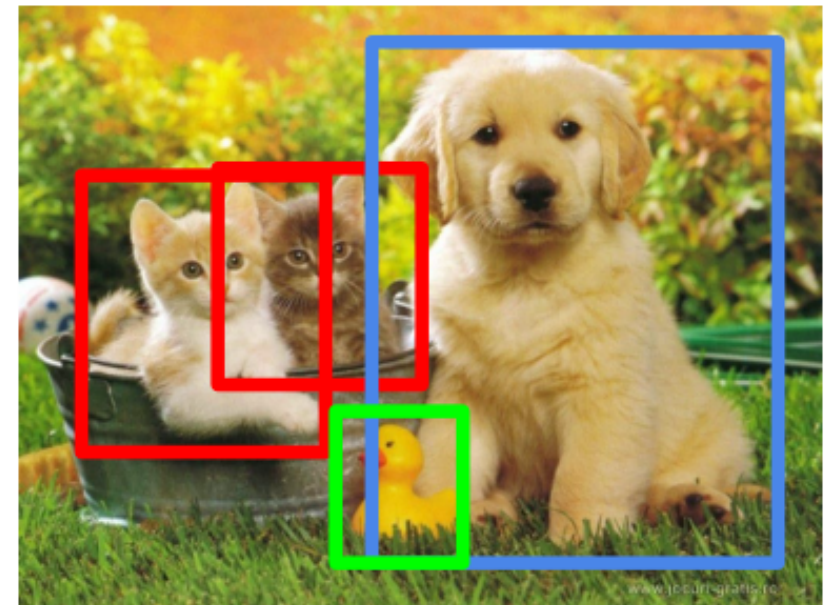
Hongkai Zhang^{1,2}, Hong Chang^{1,2}, Bingpeng Ma², Naiyan Wang³, Xilin Chen^{1,2}

¹ Institute of Computing Technology, Chinese Academy of Sciences (ICT, CAS)

² University of Chinese Academy of Sciences (UCAS), ³ TuSimple

Object Detection

- Goal
 - **Localize** all the objects in an image
 - Decide semantic **categories** of the objects



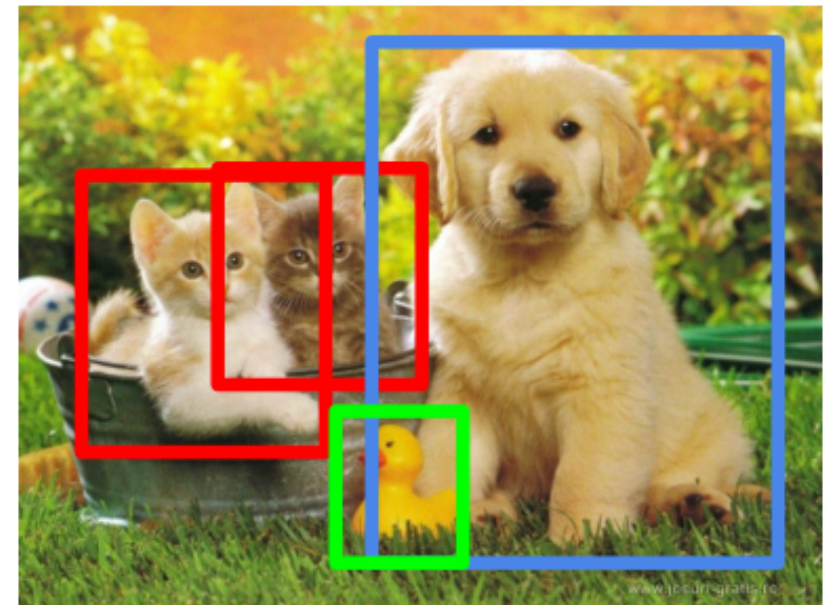
CAT, DOG, DUCK

Image source: CS231n Lecture, Stanford University.

Object detection aims to localize the objects in an image and decide their semantic categories.

Object Detection

- Two major categories
 - One-stage Detectors
 - Two-stage Detectors



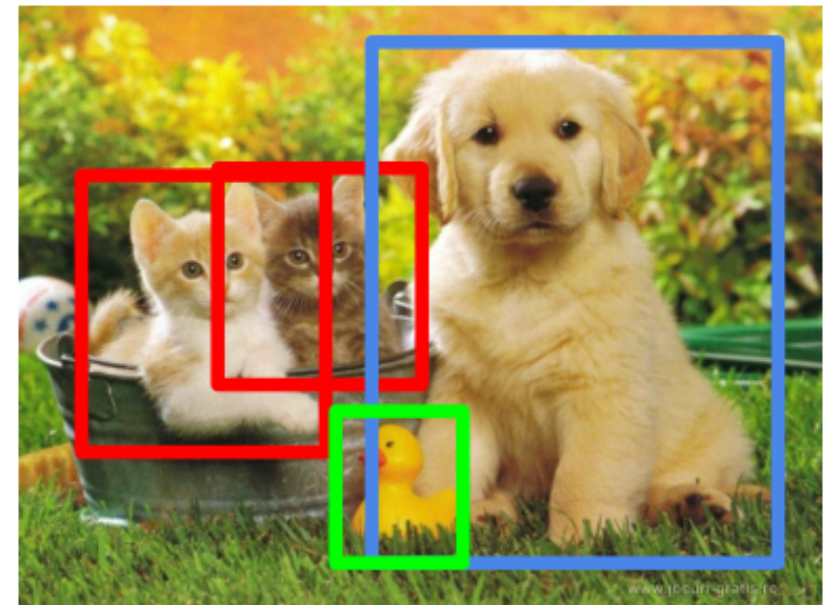
CAT, DOG, DUCK

Image source: CS231n Lecture, Stanford University.

Modern detection frameworks can be divided into two categories of one-stage detectors and two-stage detectors.

Object Detection

- Two major categories
 - One-stage Detectors
 - **Two-stage Detectors**



CAT, DOG, DUCK

Image source: CS231n Lecture, Stanford University.

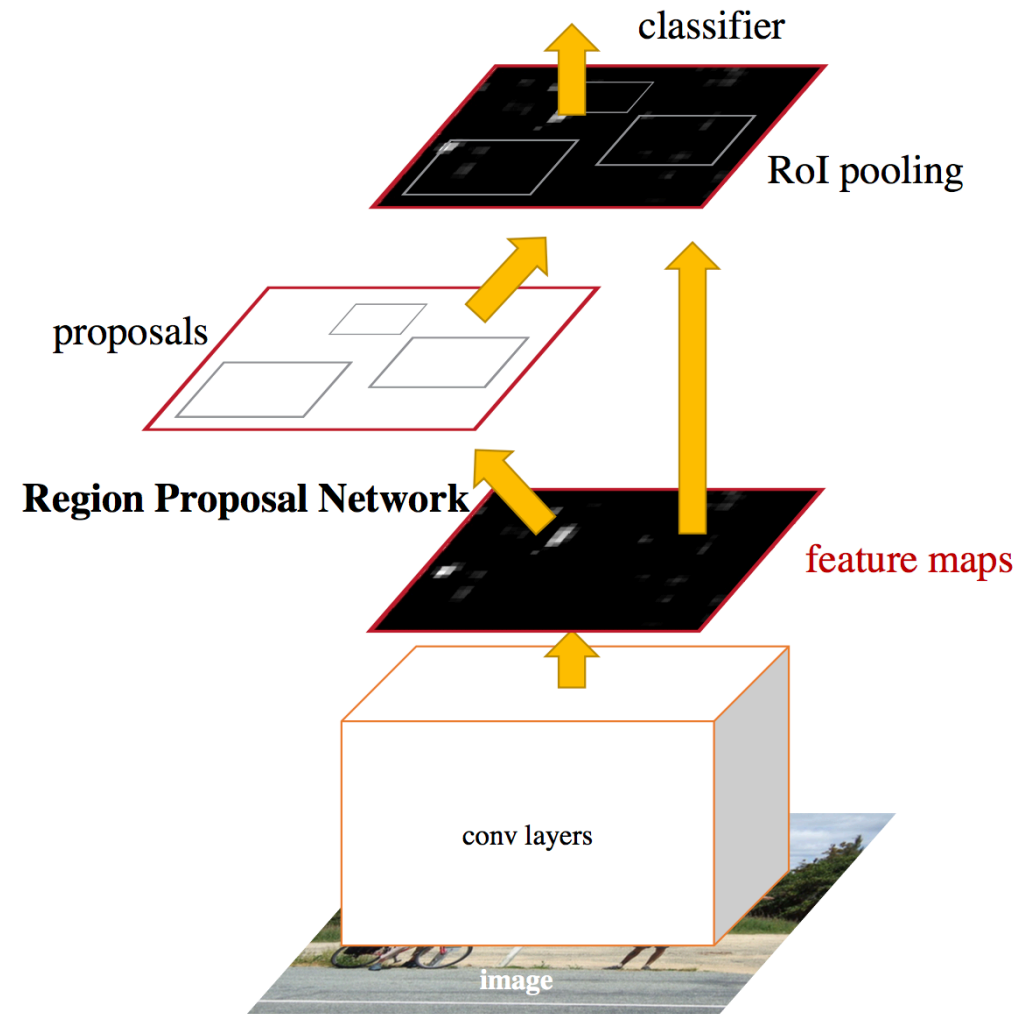
In this work, we focus on the second one.

High Quality Object Detection

- **Our goal:** High Quality Object Detection
- What is "high quality"?
 - Generally speaking, it stands for the **results under high IoU**
- Why this goal?
 - The higher the accuracy of the result -> the better

Faster R-CNN [1]

- Two-stage detectors

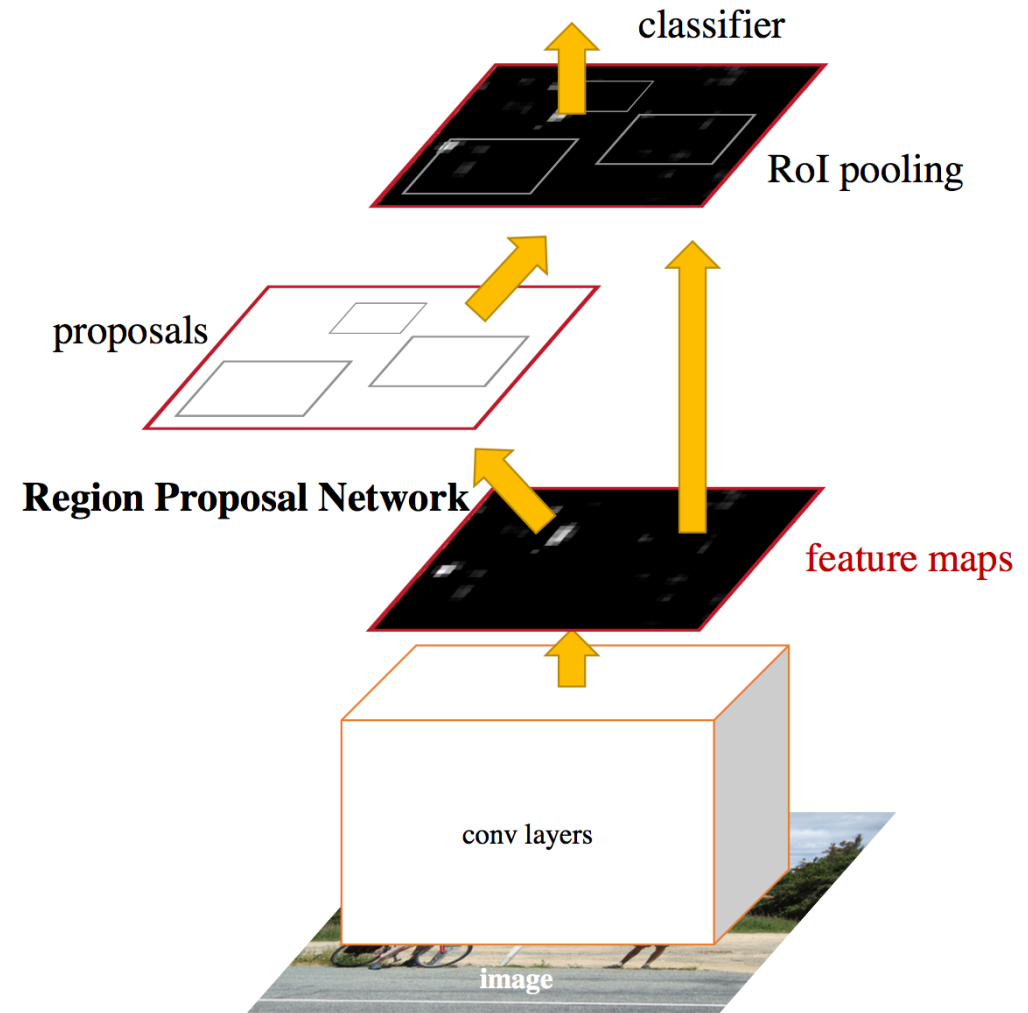


[1] Shaoqing Ren, et al. Faster R-CNN: Towards real-time object detection with region proposal networks. NIPS 2015.

To better understand our work, let's first recall the structure of a representative two-stage detector Faster R-CNN.

Faster R-CNN [1]

- Two-stage detectors
 - Coarse-to-fine manner
 - Stage 1: Anchor -> Proposal
 - Stage 2: **Proposal** -> Prediction



[1] Shaoqing Ren, et al. Faster R-CNN: Towards real-time object detection with region proposal networks. NIPS 2015.

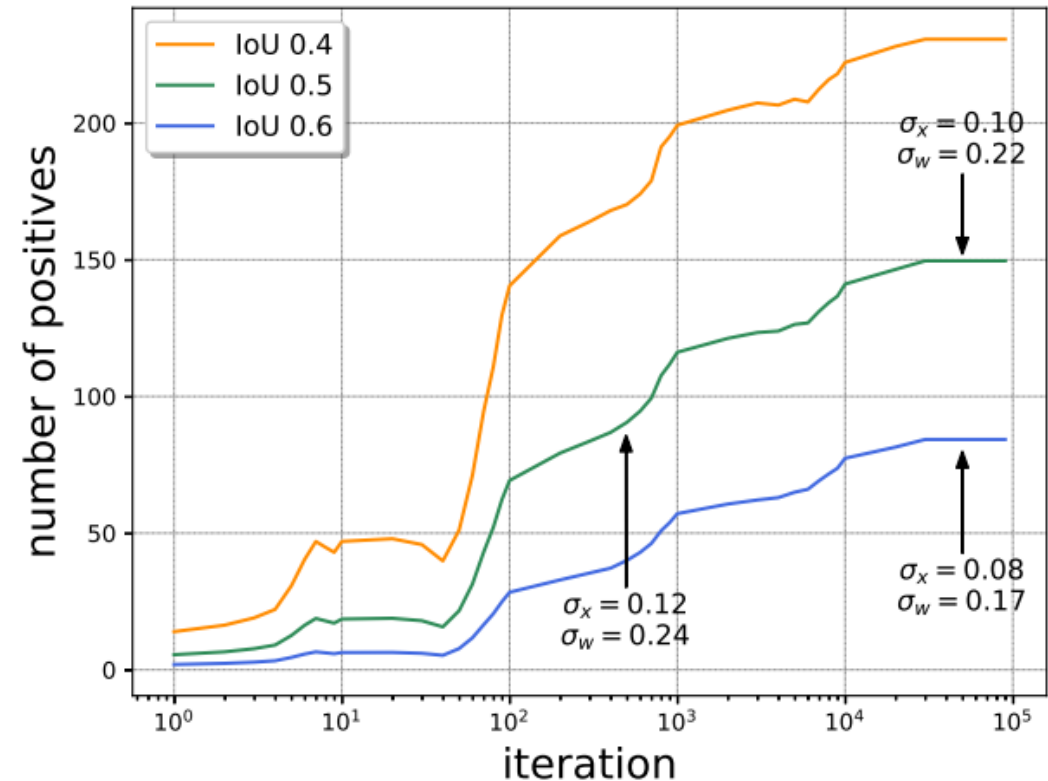
Faster R-CNN follows a coarse-to-fine manner which first refines anchors to get the region proposals, then refines proposals and get the final predictions.

Inconsistency Problem in Faster R-CNN

- The **dynamic training** and **fixed settings** are inconsistent

Inconsistency Problem in Faster R-CNN

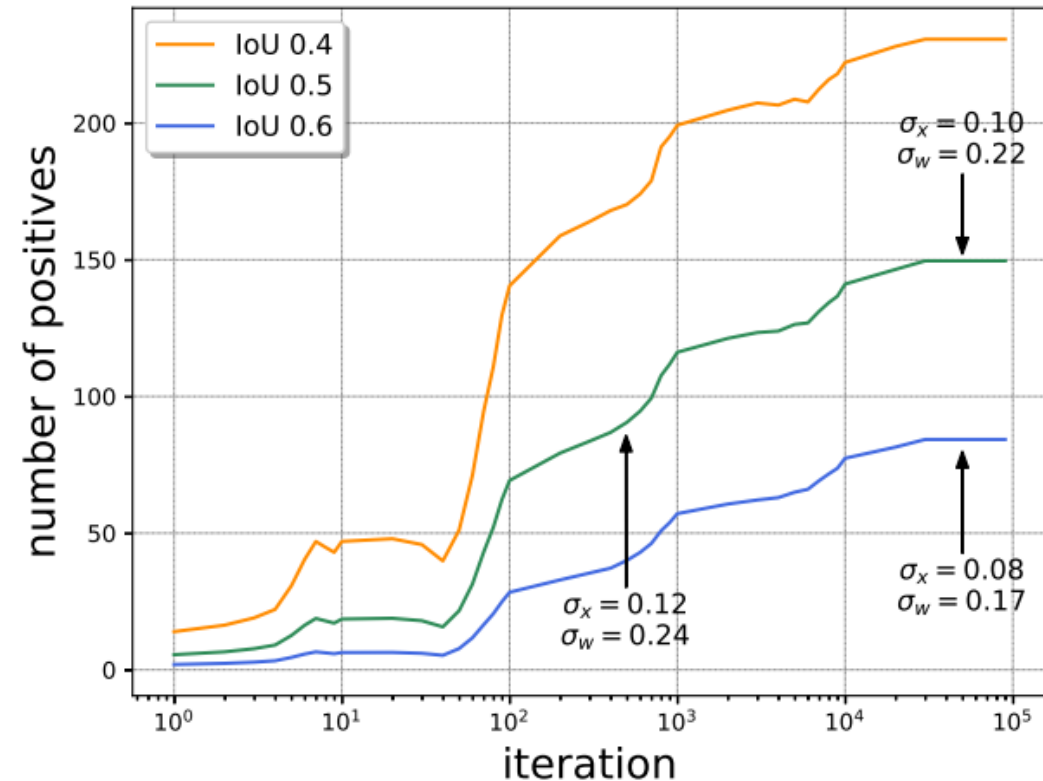
- The dynamic training and fixed settings are inconsistent
 - **Dynamic:** R-CNN input is the proposals



Specifically, the input of R-CNN is the proposals. As the training goes, the quality of proposals actually improves regardless of the certain IoU.

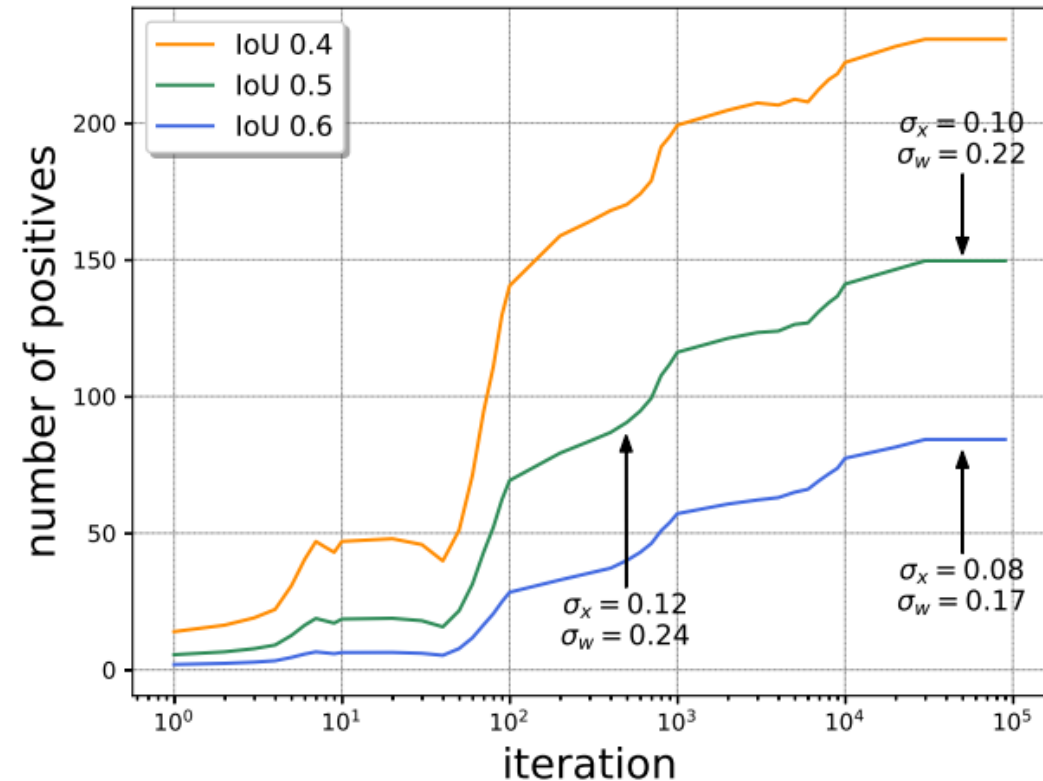
Inconsistency Problem in Faster R-CNN

- The dynamic training and fixed settings are inconsistent
 - Dynamic: R-CNN input is the proposals
 - Fixed: network settings
- Harmful to high quality Object Detection



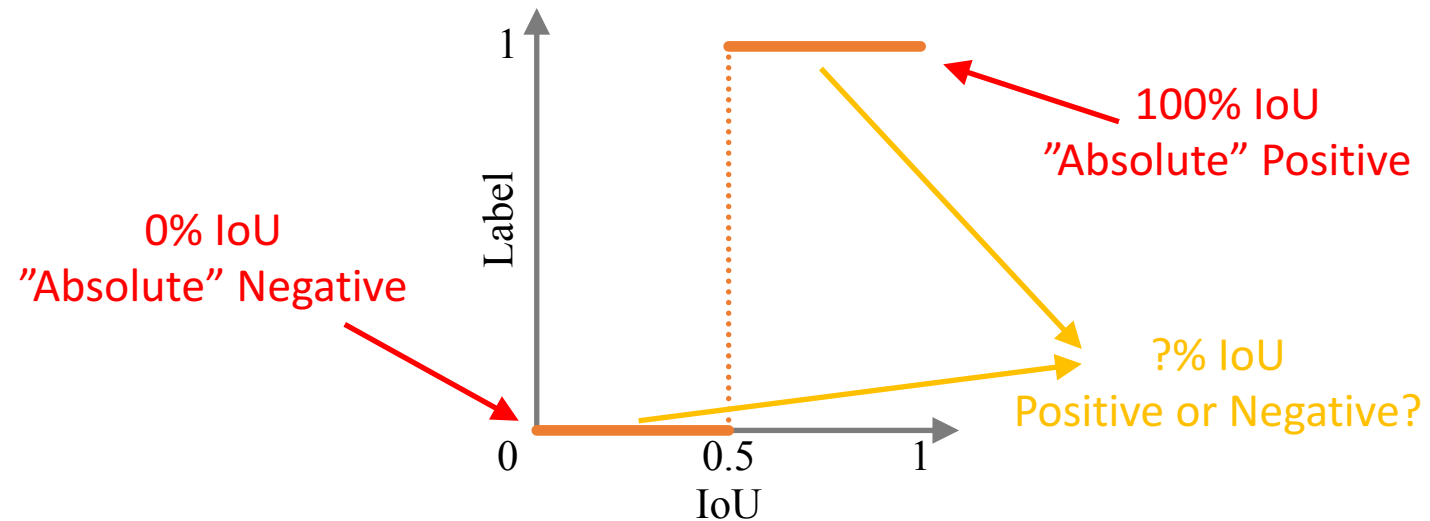
Inconsistency Problem in Faster R-CNN

- The dynamic training and fixed settings are inconsistent
 - Dynamic: R-CNN input is the proposals
 - Fixed: network settings
- Harmful to high quality Object Detection
 - Proposal classification
 - Bounding box regression



Proposal Classification

- **How to assign labels** is an interesting question
 - Assignment strategy in Faster R-CNN ($T_+ = T_- = 0.5$)



- Why the threshold is fixed at 0.5?

Proposal Classification

- Training with different **IoU thresholds** will lead to classifiers with **corresponding quality** [2]

Backbone	IoU	AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
ResNet-50-FPN	0.4	35.4	58.2	53.0	44.1	29.2	7.3
ResNet-50-FPN	0.5	36.6	58.1	53.5	45.8	31.5	8.8
ResNet-50-FPN	0.6	35.7	56.0	51.6	44.5	31.6	9.3

[2] Zhaowei Cai, et al. Cascade R-CNN: Delving into High Quality Object Detection. CVPR 2018.

As mentioned in Cascade R-CNN, training with different IoU thresholds will lead to classifiers with corresponding quality.

Proposal Classification

- High quality object detection requires **high IoU threshold**, but directly raising the IoU threshold is impractical due to the **vanishing positive samples** [2]

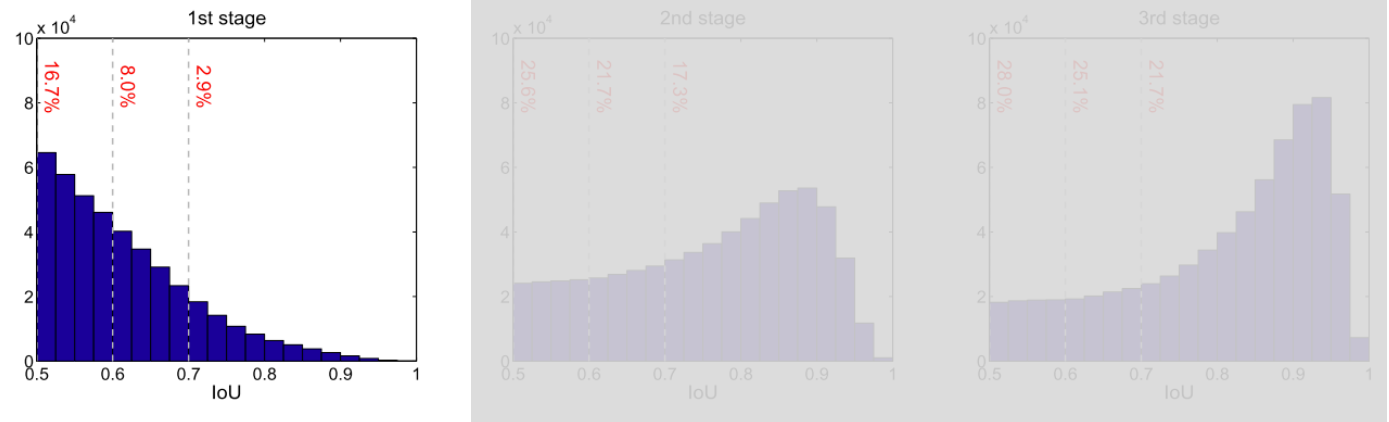


Fig. 4: IoU histograms of training samples of each cascade stage. The distribution of the 1st stage is the RPN output. Shown in red are the percentage of positives for the corresponding IoU threshold.

[2] Zhaowei Cai, et al. Cascade R-CNN: Delving into High Quality Object Detection. CVPR 2018.

To train a high quality classifier we need a high IoU threshold, but directly raising it will lead to overfitting due to the vanishing positives.

Proposal Classification

- Cascading several stages to lift the IoU of proposals [2]
- Effective yet **time-consuming**

- **Better ideas?**

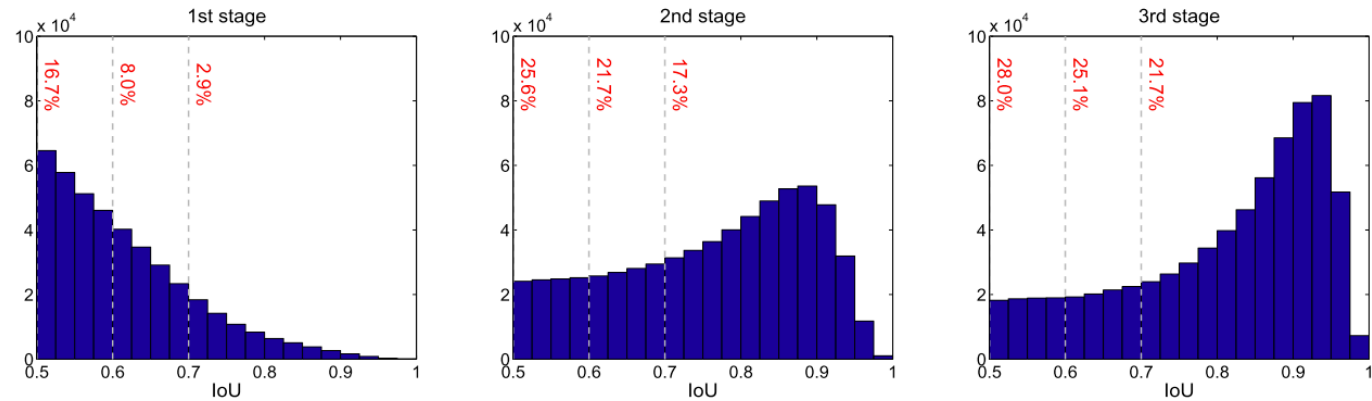


Fig. 4: IoU histograms of training samples of each cascade stage. The distribution of the 1st stage is the RPN output. Shown in red are the percentage of positives for the corresponding IoU threshold.

[2] Zhaowei Cai, et al. Cascade R-CNN: Delving into High Quality Object Detection. CVPR 2018.

So Cascade R-CNN adopts several sequential stages to lift the IoU of proposals, which is effective yet time-consuming. So are there any better ideas?

Proposal Classification

- The quality of proposals actually **improves along the training process**
- It inspires us to take a **progressive** approach in training

Proposal Classification

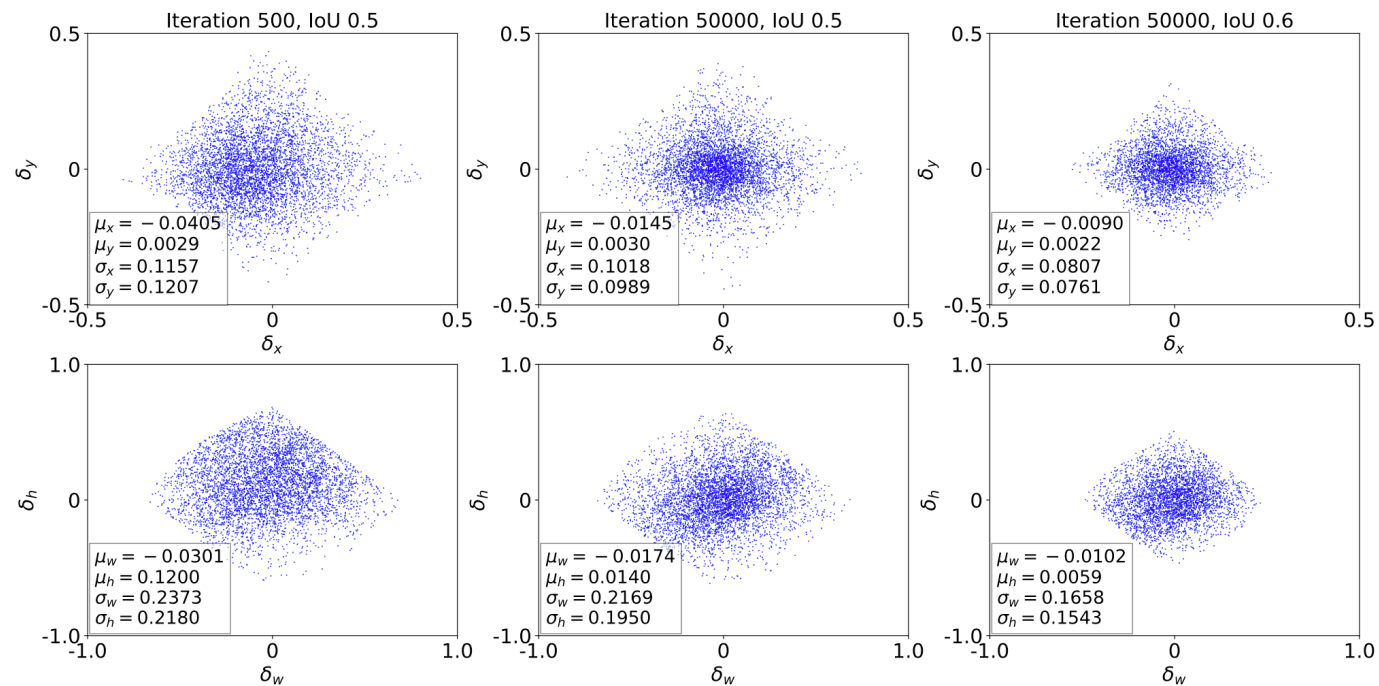
- The quality of proposals actually **improves along the training process**
- It inspires us to take a **progressive** approach in training
 - At the **beginning** -> no enough high quality proposals -> low IoU threshold

Proposal Classification

- The quality of proposals actually **improves along the training process**
- It inspires us to take a **progressive** approach in training
 - At the beginning -> no enough high quality proposals -> low IoU threshold
 - As the training goes -> **proposal quality improves** -> adjust IoU threshold

Bounding Box Regression

- Regression labels are **shifting** during training, showing the **improved proposal quality**

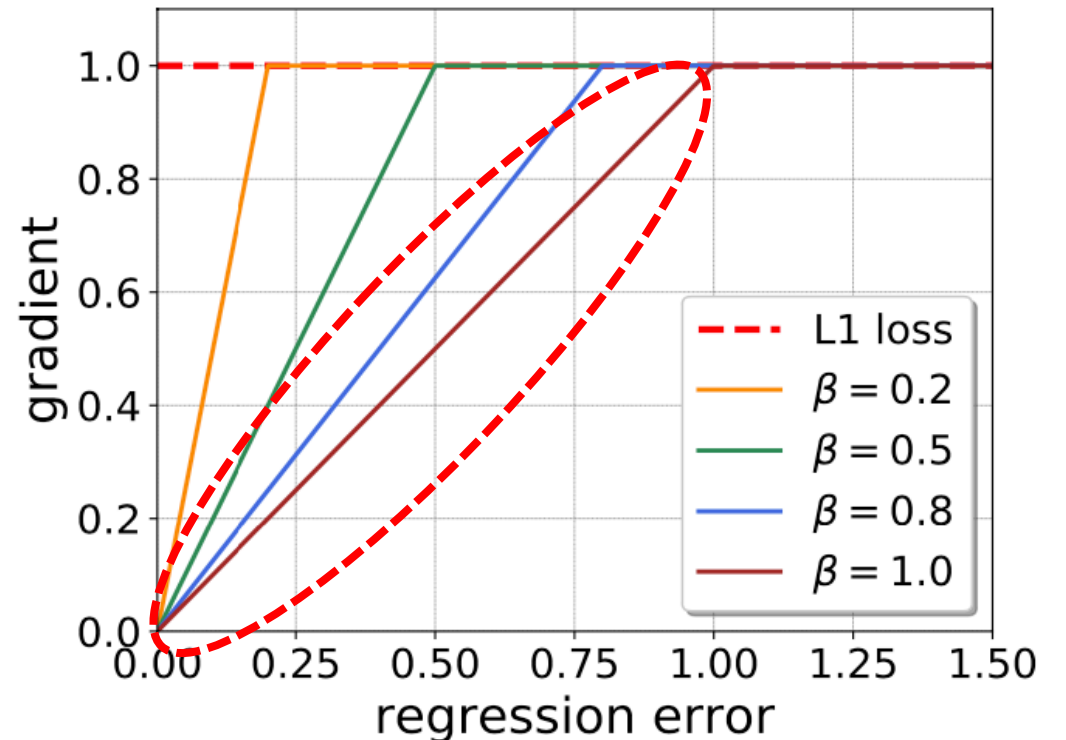


As for the regression task, we also find that the regression labels are shifting during training, showing the improved proposal quality.

Bounding Box Regression

- SmoothL1 Loss (default beta=1.0) will **reduce the contributions of high quality samples (red circle)**

$$\text{SmoothL1}(x, \beta) = \begin{cases} 0.5|x|^2/\beta, & \text{if } |x| < \beta, \\ |x| - 0.5\beta, & \text{otherwise.} \end{cases}$$

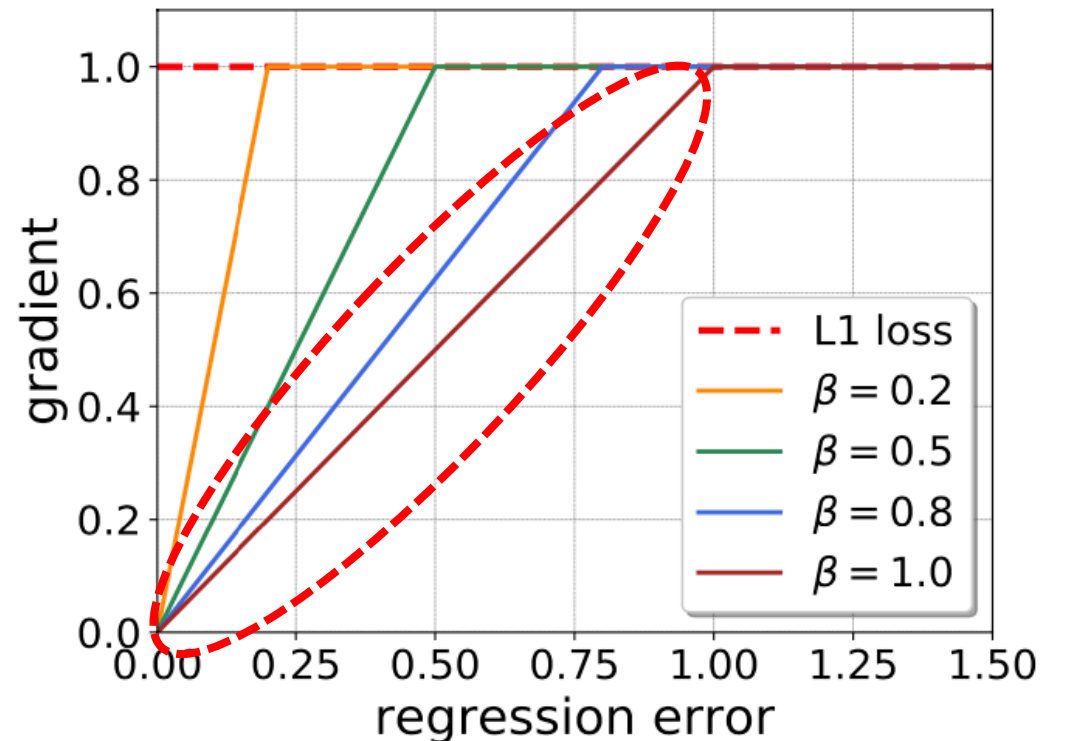


However, the regression loss function will reduce the contributions of high quality samples, which is harmful to training high quality regressors.

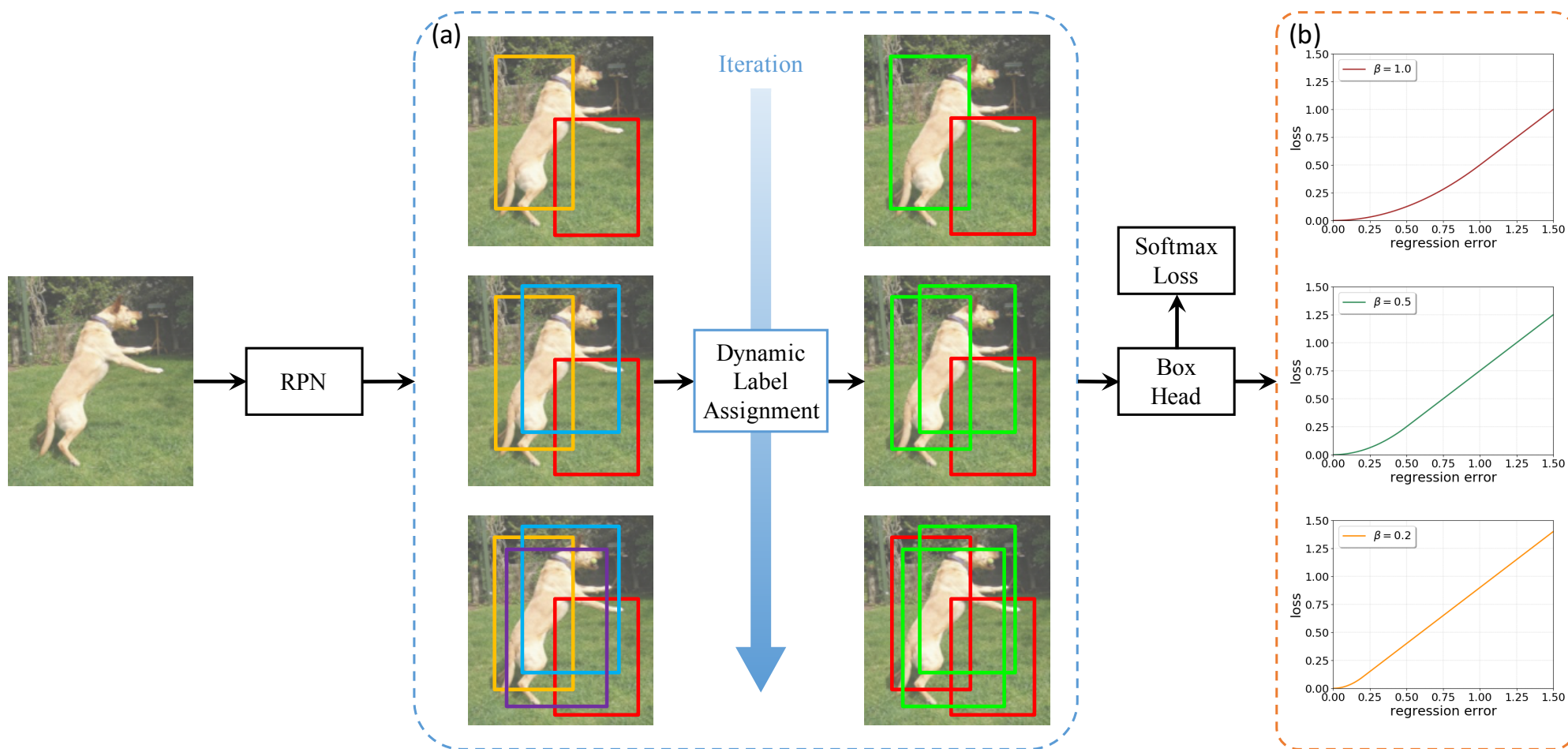
Bounding Box Regression

- SmoothL1 Loss (default beta=1.0) will reduce the contributions of high quality samples (red circle)
- **Compensate for high quality samples**

$$\text{SmoothL1}(x, \beta) = \begin{cases} 0.5|x|^2/\beta, & \text{if } |x| < \beta, \\ |x| - 0.5\beta, & \text{otherwise.} \end{cases}$$

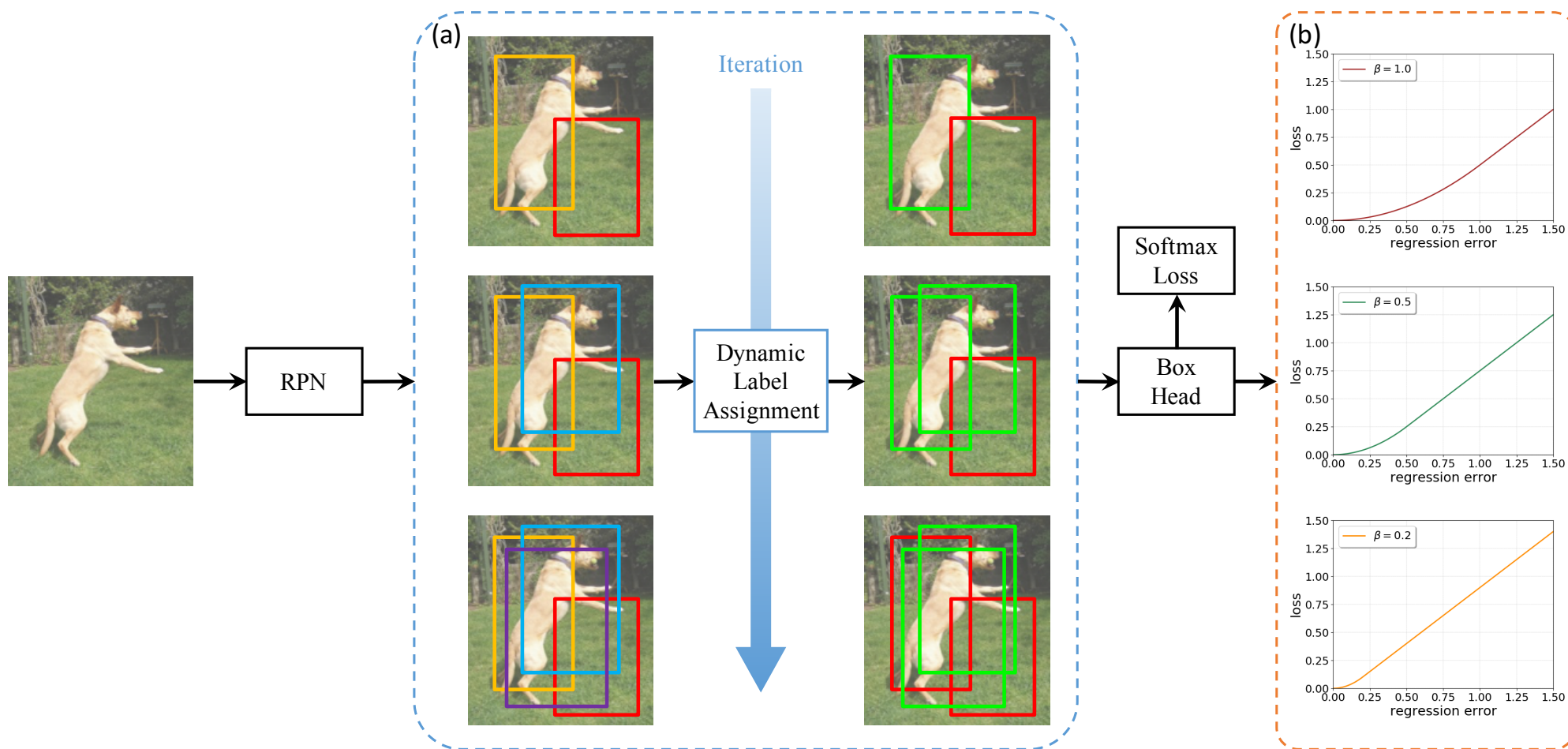


Dynamic R-CNN



To better exploit the dynamic property of the training procedure, we propose Dynamic R-CNN which consists of the following two components.

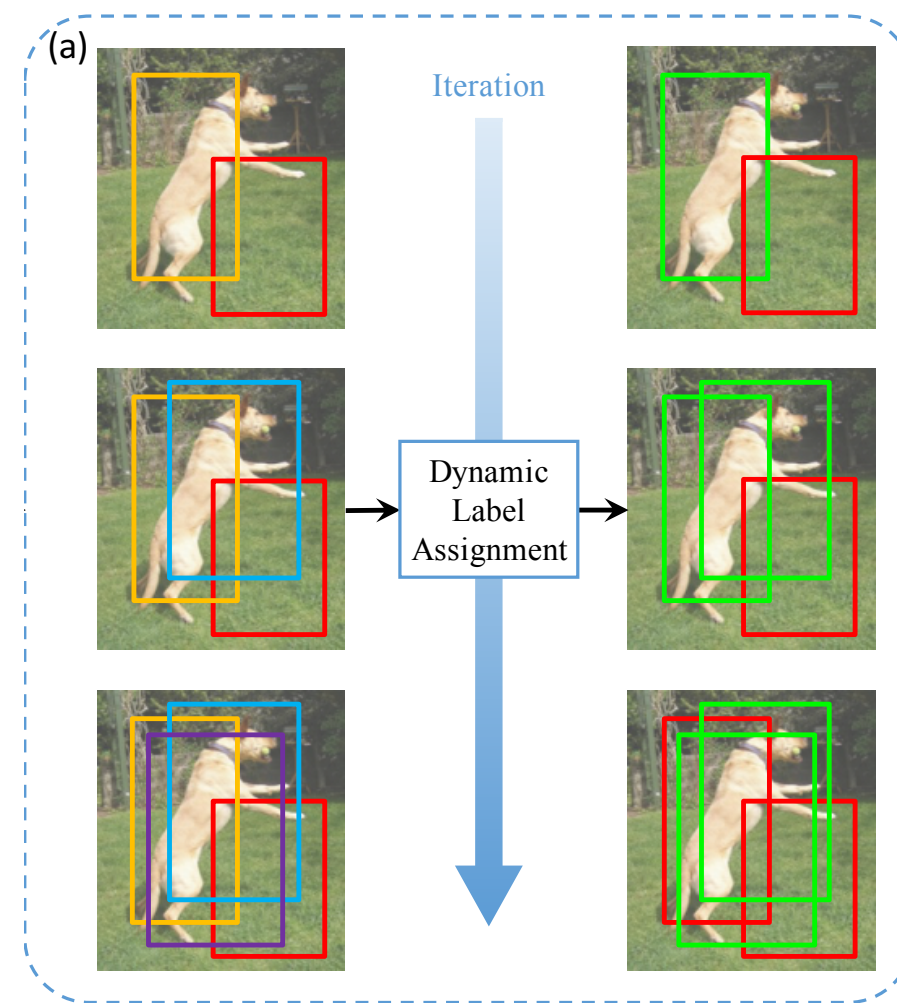
Dynamic R-CNN



Our key insight is adjusting the second stage classifier and regressor to fit the distribution change of proposals.

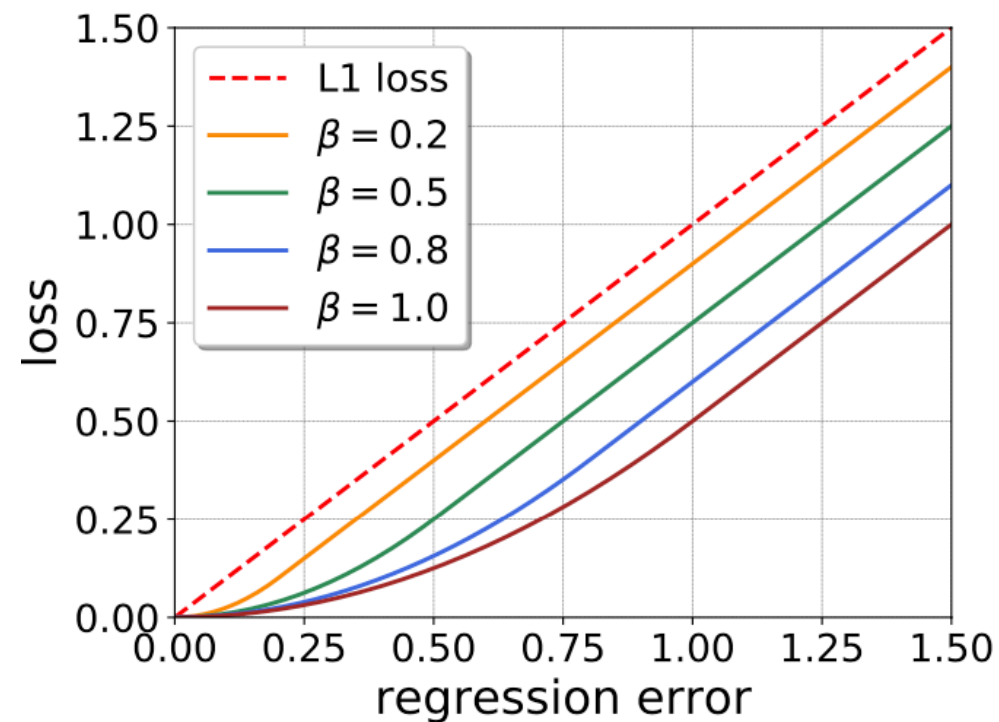
Dynamic Label Assignment (DLA) for Proposal Classification

- DLA
 - Update the training IoU threshold according to the statistics of proposals
 - Specifically, we use the K -th most accurate proposal's IoU to update the training IoU threshold



Dynamic SmoothL1 Loss (DSL) for Bounding Box Regression

- DSL
 - DSL will **update the beta of SmoothL1 Loss according to the statistics of proposals**
 - Specifically, we use the K_{β} -th most accurate proposal's regression label to update the SmoothL1 beta, then the shape of loss function will be changed



Dynamic R-CNN Algorithm

Algorithm 1 Dynamic R-CNN

Input:

Proposal set P , ground-truth set G .

IoU threshold top-k K_I , β top-k K_β , update iteration count C .

Output:

Trained object detector D .

- 1: Initialize IoU threshold and SmoothL1 β as T_{now}, β_{now}
 - 2: Build two empty sets $\mathcal{S}_I, \mathcal{S}_E$ for recording the IoUs and regression labels
 - 3: **for** $i = 0$ to max_iter **do**
 - 4: Obtain matched IoUs I and regression labels E between P and G
 - 5: Select thresholds I_k, E_k based on the K_I, K_β
 - 6: Record corresponding values, add I_k to \mathcal{S}_I and E_k to \mathcal{S}_E
 - 7: **if** $i \% C == 0$ **then**
 - 8: Update IoU threshold: $T_{now} = \text{Mean}(\mathcal{S}_I)$
 - 9: Update SmoothL1 β : $\beta_{now} = \text{Median}(\mathcal{S}_E)$
 - 10: $\mathcal{S}_I = \emptyset, \mathcal{S}_E = \emptyset$
 - 11: Train the network with T_{now}, β_{now}
 - 12: **return** Improved object detector D
-

Dynamic R-CNN Algorithm

Algorithm 1 Dynamic R-CNN

Input:

Proposal set P , ground-truth set G .

IoU threshold top-k K_I , β top-k K_β , update iteration count C .

Output:

Trained object detector D .

- 1: Initialize IoU threshold and SmoothL1 β as T_{now}, β_{now}
 - 2: Build two empty sets $\mathcal{S}_I, \mathcal{S}_E$ for recording the IoUs and regression labels
 - 3: **for** $i = 0$ to max_iter **do**
 - 4: Obtain matched IoUs I and regression labels E between P and G
 - 5: Select thresholds I_k, E_k based on the K_I, K_β
 - 6: Record corresponding values, add I_k to \mathcal{S}_I and E_k to \mathcal{S}_E
 - 7: **if** $i \% C == 0$ **then**
 - 8: Update IoU threshold: $T_{now} = \text{Mean}(\mathcal{S}_I)$
 - 9: Update SmoothL1 β : $\beta_{now} = \text{Median}(\mathcal{S}_E)$
 - 10: $\mathcal{S}_I = \emptyset, \mathcal{S}_E = \emptyset$
 - 11: Train the network with T_{now}, β_{now}
 - 12: **return** Improved object detector D
-

Dynamic R-CNN Algorithm

Algorithm 1 Dynamic R-CNN

Input:

Proposal set P , ground-truth set G .

IoU threshold top-k K_I , β top-k K_β , update iteration count C .

Output:

Trained object detector D .

- 1: Initialize IoU threshold and SmoothL1 β as T_{now}, β_{now}
 - 2: Build two empty sets $\mathcal{S}_I, \mathcal{S}_E$ for recording the IoUs and regression labels
 - 3: **for** $i = 0$ to max_iter **do**
 - 4: Obtain matched IoUs I and regression labels E between P and G
 - 5: Select thresholds I_k, E_k based on the K_I, K_β
 - 6: Record corresponding values, add I_k to \mathcal{S}_I and E_k to \mathcal{S}_E
 - 7: **if** $i \% C == 0$ **then**
 - 8: Update IoU threshold: $T_{now} = \text{Mean}(\mathcal{S}_I)$
 - 9: Update SmoothL1 β : $\beta_{now} = \text{Median}(\mathcal{S}_E)$
 - 10: $\mathcal{S}_I = \emptyset, \mathcal{S}_E = \emptyset$
 - 11: Train the network with T_{now}, β_{now}
 - 12: **return** Improved object detector D
-

Dynamic R-CNN Algorithm

Algorithm 1 Dynamic R-CNN

Input:

Proposal set P , ground-truth set G .

IoU threshold top-k K_I , β top-k K_β , update iteration count C .

Output:

Trained object detector D .

- 1: Initialize IoU threshold and SmoothL1 β as T_{now}, β_{now}
 - 2: Build two empty sets $\mathcal{S}_I, \mathcal{S}_E$ for recording the IoUs and regression labels
 - 3: **for** $i = 0$ to max_iter **do**
 - 4: Obtain matched IoUs I and regression labels E between P and G
 - 5: Select thresholds I_k, E_k based on the K_I, K_β
 - 6: Record corresponding values, add I_k to \mathcal{S}_I and E_k to \mathcal{S}_E
 - 7: **if** $i \% C == 0$ **then**
 - 8: Update IoU threshold: $T_{now} = \text{Mean}(\mathcal{S}_I)$
 - 9: Update SmoothL1 β : $\beta_{now} = \text{Median}(\mathcal{S}_E)$
 - 10: $\mathcal{S}_I = \emptyset, \mathcal{S}_E = \emptyset$
 - 11: Train the network with T_{now}, β_{now}
 - 12: **return** Improved object detector D
-

Then, Dynamic R-CNN will update the IoU threshold and SmoothL1 beta every C iterations to fit the dynamic property during training.

Dynamic R-CNN Algorithm

Algorithm 1 Dynamic R-CNN

Input:

Proposal set P , ground-truth set G .

IoU threshold top-k K_I , β top-k K_β , update iteration count C .

Output:

Trained object detector D .

- 1: Initialize IoU threshold and SmoothL1 β as T_{now}, β_{now}
 - 2: Build two empty sets $\mathcal{S}_I, \mathcal{S}_E$ for recording the IoUs and regression labels
 - 3: **for** $i = 0$ to max_iter **do**
 - 4: Obtain matched IoUs I and regression labels E between P and G
 - 5: Select thresholds I_k, E_k based on the K_I, K_β
 - 6: Record corresponding values, add I_k to \mathcal{S}_I and E_k to \mathcal{S}_E
 - 7: **if** $i \% C == 0$ **then**
 - 8: Update IoU threshold: $T_{now} = \text{Mean}(\mathcal{S}_I)$
 - 9: Update SmoothL1 β : $\beta_{now} = \text{Median}(\mathcal{S}_E)$
 - 10: $\mathcal{S}_I = \emptyset, \mathcal{S}_E = \emptyset$
 - 11: Train the network with T_{now}, β_{now}
 - 12: **return** Improved object detector D
-

Finally, we get an improved object detector.

Experiments (Main Result)

Table 1. Comparisons with different baselines (our re-implementations) on COCO test-dev set. “MST” and “*” stand for multi-scale training and testing respectively. “2×” and “3×” are training schedules which extend the iterations by 2/3 times.

Method	Backbone	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Faster R-CNN	ResNet-50	37.3	58.5	40.6	20.3	39.2	49.1
Faster R-CNN+2×	ResNet-50	38.1	58.9	41.5	20.5	40.0	50.0
Faster R-CNN	ResNet-101	39.3	60.5	42.7	21.3	41.8	51.7
Faster R-CNN+2×	ResNet-101	39.9	60.6	43.5	21.4	42.4	52.1
Faster R-CNN+3×+MST	ResNet-101	42.8	63.8	46.8	24.8	45.6	55.6
Faster R-CNN+3×+MST	ResNet-101-DCN	44.8	65.5	48.8	26.2	47.6	58.1
Faster R-CNN+3×+MST*	ResNet-101-DCN	46.9	68.1	51.4	30.6	49.6	58.1
Dynamic R-CNN	ResNet-50	39.1	58.0	42.8	21.3	40.9	50.3
Dynamic R-CNN+2×	ResNet-50	39.9	58.6	43.7	21.6	41.5	51.9
Dynamic R-CNN	ResNet-101	41.2	60.1	45.1	22.5	43.6	53.2
Dynamic R-CNN+2×	ResNet-101	42.0	60.7	45.9	22.7	44.3	54.3
Dynamic R-CNN+3×+MST	ResNet-101	44.7	63.6	49.1	26.0	47.4	57.2
Dynamic R-CNN+3×+MST	ResNet-101-DCN	46.9	65.9	51.3	28.1	49.6	60.0
Dynamic R-CNN+3×+MST*	ResNet-101-DCN	49.2	68.6	54.0	32.5	51.7	60.3

Experiments (Main Result)

Table 1. Comparisons with different baselines (our re-implementations) on COCO test-dev set. “MST” and “*” stand for multi-scale training and testing respectively. “2×” and “3×” are training schedules which extend the iterations by 2/3 times.

Method	Backbone	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Faster R-CNN	ResNet-50	37.3	58.5	40.6	20.3	39.2	49.1
Faster R-CNN+2×	ResNet-50	38.1	58.9	41.5	20.5	40.0	50.0
Faster R-CNN	ResNet-101	39.3	60.5	42.7	21.3	41.8	51.7
Faster R-CNN+2×	ResNet-101	39.9	60.6	43.5	21.4	42.4	52.1
Faster R-CNN+3×+MST	ResNet-101	42.8	63.8	46.8	24.8	45.6	55.6
Faster R-CNN+3×+MST	ResNet-101-DCN	44.8	65.5	48.8	26.2	47.6	58.1
Faster R-CNN+3×+MST*	ResNet-101-DCN	46.9	68.1	51.4	30.6	49.6	58.1
Dynamic R-CNN	ResNet-50	39.1	58.0	42.8	21.3	40.9	50.3
Dynamic R-CNN+2×	ResNet-50	39.9	58.6	43.7	21.6	41.5	51.9
Dynamic R-CNN	ResNet-101	41.2	60.1	45.1	22.5	43.6	53.2
Dynamic R-CNN+2×	ResNet-101	42.0	60.7	45.9	22.7	44.3	54.3
Dynamic R-CNN+3×+MST	ResNet-101	44.7	63.6	49.1	26.0	47.4	57.2
Dynamic R-CNN+3×+MST	ResNet-101-DCN	46.9	65.9	51.3	28.1	49.6	60.0
Dynamic R-CNN+3×+MST*	ResNet-101-DCN	49.2	68.6	54.0	32.5	51.7	60.3

As shown in Table 1, our method can work on different backbones and it is also compatible with other training and testing skills.

Experiments (Main Result)

Table 1. Comparisons with different baselines (our re-implementations) on COCO test-dev set. “MST” and “*” stand for multi-scale training and testing respectively. “2×” and “3×” are training schedules which extend the iterations by 2/3 times.

Method	Backbone	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Faster R-CNN	ResNet-50	37.3	58.5	40.6	20.3	39.2	49.1
Faster R-CNN+2×	ResNet-50	38.1	58.9	41.5	20.5	40.0	50.0
Faster R-CNN	ResNet-101	39.3	60.5	42.7	21.3	41.8	51.7
Faster R-CNN+2×	ResNet-101	39.9	60.6	43.5	21.4	42.4	52.1
Faster R-CNN+3×+MST	ResNet-101	42.8	63.8	46.8	24.8	45.6	55.6
Faster R-CNN+3×+MST	ResNet-101-DCN	44.8	65.5	48.8	26.2	47.6	58.1
Faster R-CNN+3×+MST*	ResNet-101-DCN	46.9	68.1	51.4	30.6	49.6	58.1
Dynamic R-CNN	ResNet-50	39.1	58.0	42.8	21.3	40.9	50.3
Dynamic R-CNN+2×	ResNet-50	39.9	58.6	43.7	21.6	41.5	51.9
Dynamic R-CNN	ResNet-101	41.2	60.1	45.1	22.5	43.6	53.2
Dynamic R-CNN+2×	ResNet-101	42.0	60.7	45.9	22.7	44.3	54.3
Dynamic R-CNN+3×+MST	ResNet-101	44.7	63.6	49.1	26.0	47.4	57.2
Dynamic R-CNN+3×+MST	ResNet-101-DCN	46.9	65.9	51.3	28.1	49.6	60.0
Dynamic R-CNN+3×+MST*	ResNet-101-DCN	49.2	68.6	54.0	32.5	51.7	60.3

Generally speaking, our method can improve different baselines by almost 2 points AP consistently.

Experiments (Components)

Table 2. Results of each component in Dynamic R-CNN on COCO val set.

Backbone	DLA	DSL	AP	Δ AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
ResNet-50-FPN			37.0	-	58.0	53.5	46.0	32.6	9.7
ResNet-50-FPN		✓	38.0	+1.0	57.6	53.5	46.7	34.4	13.2
ResNet-50-FPN	✓		38.2	+1.2	57.5	53.6	47.1	35.2	12.6
ResNet-50-FPN	✓	✓	38.9	+1.9	57.3	53.6	47.4	36.3	15.2

Then, we conduct experiments to validate the effectiveness of different components in Dynamic R-CNN.

Experiments (Components)

Table 2. Results of each component in Dynamic R-CNN on COCO val set.

Backbone	DLA	DSL	AP	Δ AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
ResNet-50-FPN			37.0	-	58.0	53.5	46.0	32.6	9.7
ResNet-50-FPN		✓	38.0	+1.0	57.6	53.5	46.7	34.4	13.2
ResNet-50-FPN	✓		38.2	+1.2	57.5	53.6	47.1	35.2	12.6
ResNet-50-FPN	✓	✓	38.9	+1.9	57.3	53.6	47.4	36.3	15.2

Experiments (Components)

Table 2. Results of each component in Dynamic R-CNN on COCO val set.

Backbone	DLA	DSL	AP	Δ AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
ResNet-50-FPN			37.0	-	58.0	53.5	46.0	32.6	9.7
ResNet-50-FPN		✓	38.0	+1.0	57.6	53.5	46.7	34.4	13.2
ResNet-50-FPN	✓		38.2	+1.2	57.5	53.6	47.1	35.2	12.6
ResNet-50-FPN	✓	✓	38.9	+1.9	57.3	53.6	47.4	36.3	15.2

Since with DLA the contributions of positives are reduced even more, DLA and DSL can work together.

Experiments (Components)

Table 2. Results of each component in Dynamic R-CNN on COCO val set.

Backbone	DLA	DSL	AP	Δ AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
ResNet-50-FPN			37.0	-	58.0	53.5	46.0	32.6	9.7
ResNet-50-FPN		✓	38.0	+1.0	57.6	53.5	46.7	34.4	13.2
ResNet-50-FPN	✓		38.2	+1.2	57.5	53.6	47.1	35.2	12.6
ResNet-50-FPN	✓	✓	38.9	+1.9	57.3	53.6	47.4	36.3	15.2

+5.5
AP90

To sum up, Dynamic R-CNN improves the baseline by 1.9 points AP and 5.5 points AP90.

Experiments (Dynamic trends)

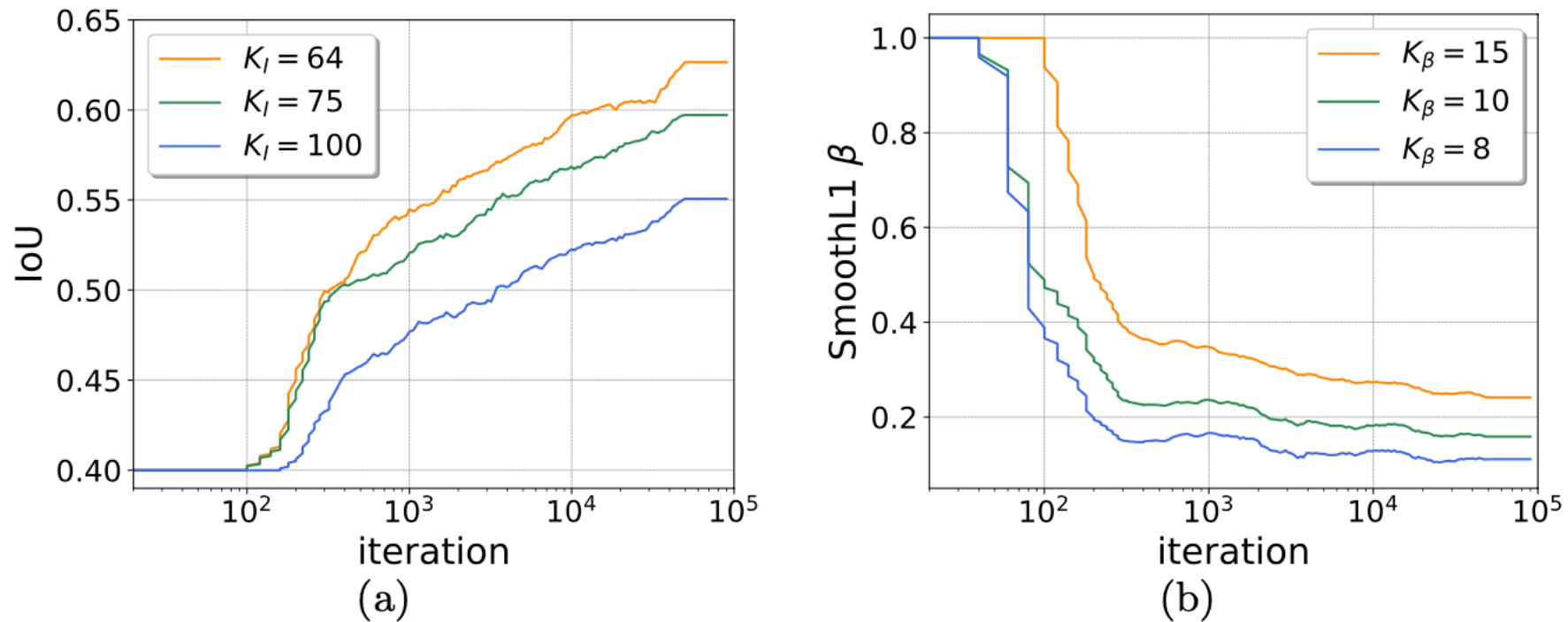


Fig. 5. Trends of (a) IoU threshold and (b) SmoothL1 β under different settings based on our method. Obviously the distribution has changed a lot during training.

Experiments (Dynamic trends)

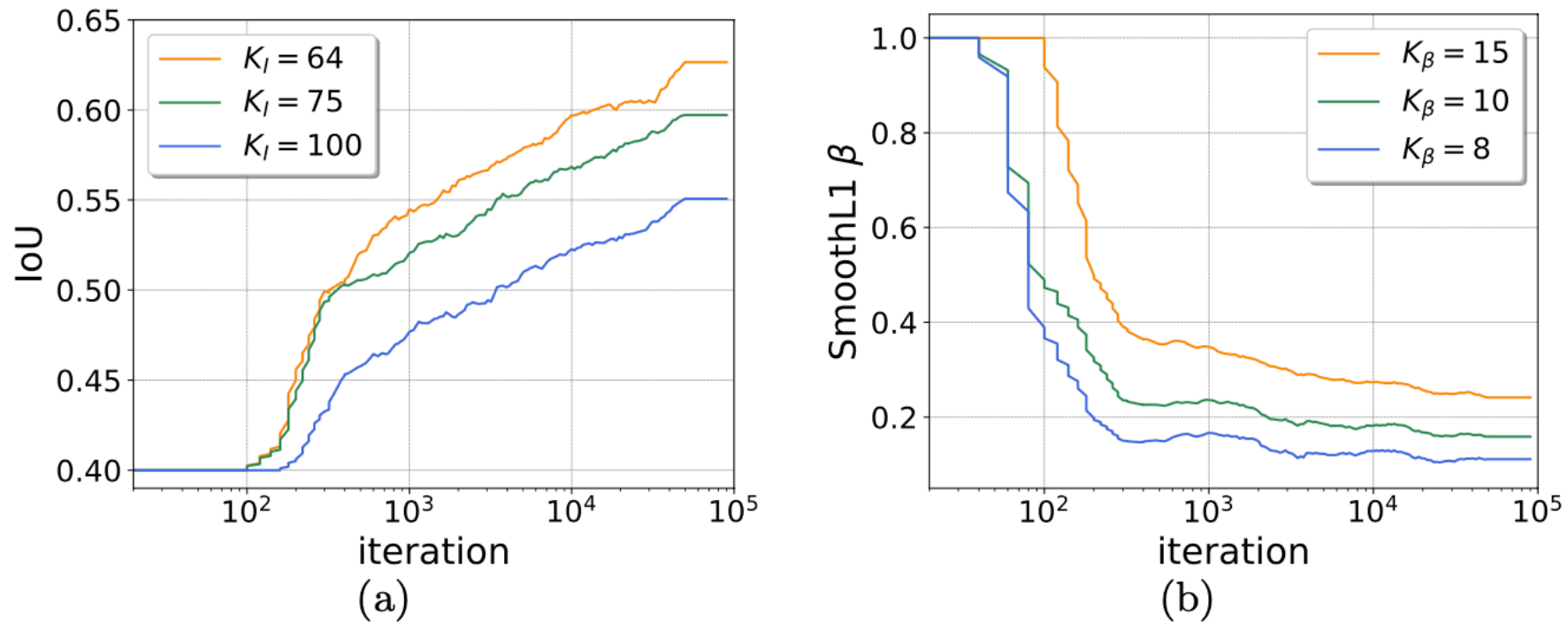


Fig. 5. Trends of (a) IoU threshold and (b) SmoothL1 β under different settings based on our method. Obviously the distribution has changed a lot during training.

Regardless of the specific settings, the trend of IoU threshold is increasing while that for beta is decreasing during training as expected.

Experiments (Hyper-parameters)

Table 3. Ablation study on K_I .

K_I	AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
-	37.0	58.0	53.5	46.0	32.6	9.7
64	38.1	57.2	53.3	46.8	35.1	12.8
75	38.2	57.5	53.6	47.1	35.2	12.6
100	37.9	57.9	53.8	46.9	34.2	11.6

Table 4. Ablation study on C .

C	AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
-	37.0	58.0	53.5	46.0	32.6	9.7
20	38.0	57.4	53.5	47.0	35.0	12.5
100	38.2	57.5	53.6	47.1	35.2	12.6
500	38.1	57.6	53.5	47.2	34.8	12.6

Table 5. Ablation study on K_β .

Setting	AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
$\beta = 1.0$	37.0	58.0	53.5	46.0	32.6	9.7
$\beta = 2.0$	35.9	57.7	53.2	45.1	30.1	8.3
$\beta = 0.5$	37.5	57.6	53.3	46.4	33.5	11.3
$K_\beta = 15$	37.6	57.3	53.1	46.0	33.9	12.5
$K_\beta = 10$	38.0	57.6	53.5	46.7	34.4	13.2
$K_\beta = 8$	37.6	57.5	53.3	45.9	33.9	12.4

Then, experiments on the effect of hyper-parameters are shown in these tables. Generally speaking, the results are not sensitive to these hyper-parameters.

Experiments (Hyper-parameters)

Table 3. Ablation study on K_I .

K_I	AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
-	37.0	58.0	53.5	46.0	32.6	9.7
64	38.1	57.2	53.3	46.8	35.1	12.8
75	38.2	57.5	53.6	47.1	35.2	12.6
100	37.9	57.9	53.8	46.9	34.2	11.6

Table 5. Ablation study on K_β .

Setting	AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
$\beta = 1.0$	37.0	58.0	53.5	46.0	32.6	9.7
$\beta = 2.0$	35.9	57.7	53.2	45.1	30.1	8.3
$\beta = 0.5$	37.5	57.6	53.3	46.4	33.5	11.3
$K_\beta = 15$	37.6	57.3	53.1	46.0	33.9	12.5
$K_\beta = 10$	38.0	57.6	53.5	46.7	34.4	13.2
$K_\beta = 8$	37.6	57.5	53.3	45.9	33.9	12.4

Experiments (Hyper-parameters)

Table 4. Ablation study on C .

C	AP	AP ₅₀	AP ₆₀	AP ₇₀	AP ₈₀	AP ₉₀
-	37.0	58.0	53.5	46.0	32.6	9.7
20	38.0	57.4	53.5	47.0	35.0	12.5
100	38.2	57.5	53.6	47.1	35.2	12.6
500	38.1	57.6	53.5	47.2	34.8	12.6

Experiments (Hyper-parameters)

- Only one additional hyper-parameter
 - Faster R-CNN IoU threshold -> KI
 - Faster R-CNN SmoothL1 beta -> Kbeta
 - Iteration count C (**additional but robust**)

Experiments (Complexity and Speed)

- Comparison to High Quality detectors (e.g. Cascade R-CNN)
 - Advantages
 - Does not increase the training time (Basically)

Experiments (Complexity and Speed)

- Comparison to High Quality detectors

(e.g. Cascade R-CNN)

- Advantages
 - Does not increase the training time (Basically)
 - **Faster inference speed**

Table 6. Inference speed comparisons using ResNet-50-FPN backbone on RTX 2080TI GPU.

Method	FPS
Dynamic R-CNN	13.9
Cascade R-CNN	11.2
Dynamic Mask R-CNN	11.3
Cascade Mask R-CNN	7.3

Experiments (Complexity and Speed)

- Comparison to High Quality detectors (e.g. Cascade R-CNN)
 - Advantages
 - Does not increase the training time (Basically)
 - Faster inference speed (**1.74 times faster** using ResNet-18 with mask head)

Table 6. Inference speed comparisons using ResNet-50-FPN backbone on RTX 2080TI GPU.

Method	FPS
Dynamic R-CNN	13.9
Cascade R-CNN	11.2
Dynamic Mask R-CNN	11.3
Cascade Mask R-CNN	7.3

Experiments (Universality)

Table 7. The universality of Dynamic R-CNN. We apply the idea of dynamic training on Mask R-CNN under different backbones. “bbox” and “segm” stand for object detection and instance segmentation results on COCO val set, respectively.

Backbone	+Dynamic	AP^{bbox}	AP_{50}^{bbox}	AP_{75}^{bbox}	AP^{segm}	AP_{50}^{segm}	AP_{75}^{segm}
ResNet-50-FPN		37.5	58.0	40.7	33.8	54.6	36.0
	✓	39.4	57.6	43.3	34.8	55.0	37.5
ResNet-101-FPN		39.7	60.7	43.2	35.6	56.9	37.7
	✓	41.8	60.4	45.8	36.7	57.5	39.4

Moreover, since our dynamic viewpoint is a general concept, we believe it can be adopted in different methods.

Experiments (Universality)

Table 7. The universality of Dynamic R-CNN. We apply the idea of dynamic training on Mask R-CNN under different backbones. “bbox” and “segm” stand for object detection and instance segmentation results on COCO val set, respectively.

Backbone	+Dynamic	AP^{bbox}	AP_{50}^{bbox}	AP_{75}^{bbox}	AP^{segm}	AP_{50}^{segm}	AP_{75}^{segm}
ResNet-50-FPN		37.5	58.0	40.7	33.8	54.6	36.0
	✓	39.4	57.6	43.3	34.8	55.0	37.5
ResNet-101-FPN		39.7	60.7	43.2	35.6	56.9	37.7
	✓	41.8	60.4	45.8	36.7	57.5	39.4

Experiments (Universality)

Table 7. The universality of Dynamic R-CNN. We apply the idea of dynamic training on Mask R-CNN under different backbones. “bbox” and “segm” stand for object detection and instance segmentation results on COCO val set, respectively.

Backbone	+Dynamic	AP^{bbox}	AP_{50}^{bbox}	AP_{75}^{bbox}	AP^{segm}	AP_{50}^{segm}	AP_{75}^{segm}
ResNet-50-FPN		37.5	58.0	40.7	33.8	54.6	36.0
	✓	39.4	57.6	43.3	34.8	55.0	37.5
ResNet-101-FPN		39.7	60.7	43.2	35.6	56.9	37.7
	✓	41.8	60.4	45.8	36.7	57.5	39.4

Experiments (State-of-the-Arts)

Table 8. Comparisons of single-model results on COCO test-dev set.

Method	Backbone	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
RetinaNet [28]	ResNet-101	39.1	59.1	42.3	21.8	42.7	50.2
CornerNet [23]	Hourglass-104	40.5	56.5	43.1	19.4	42.7	53.9
FCOS [42]	ResNet-101	41.0	60.7	44.1	24.0	44.1	51.0
FreeAnchor [49]	ResNet-101	41.8	61.1	44.9	22.6	44.7	53.9
RepPoints [46]	ResNet-101-DCN	45.0	66.1	49.0	26.6	48.6	57.5
CenterNet [50]	Hourglass-104	45.1	63.9	49.3	26.6	47.1	57.7
ATSS [48]	ResNet-101-DCN	46.3	64.7	50.4	27.7	49.8	58.4
Faster R-CNN [27]	ResNet-101	36.2	59.1	39.0	18.2	39.0	48.2
Mask R-CNN [14]	ResNet-101	38.2	60.3	41.7	20.1	41.1	50.2
Regionlets [45]	ResNet-101	39.3	59.8	-	21.7	43.7	50.9
Libra R-CNN [33]	ResNet-101	41.1	62.1	44.7	23.4	43.7	52.5
Cascade R-CNN [3]	ResNet-101	42.8	62.1	46.3	23.7	45.5	55.2
SNIP [40]	ResNet-101-DCN	44.4	66.2	49.9	27.3	47.4	56.9
DCNv2 [51]	ResNet-101-DCN	46.0	67.9	50.8	27.8	49.1	59.5
TridentNet [25]	ResNet-101-DCN	48.4	69.7	53.5	31.8	51.3	60.3
Dynamic R-CNN	ResNet-101	42.0	60.7	45.9	22.7	44.3	54.3
Dynamic R-CNN*	ResNet-101-DCN	50.1	68.3	55.6	32.8	53.0	61.2

Finally, we compare Dynamic R-CNN with the state-of-the-arts and find that our method outperforms other previous detectors.

Conclusion

- Dynamic R-CNN can bring **consistent gains with no extra overhead**, which is a **free lunch** for high quality object detection ~

Conclusion

- Dynamic R-CNN can bring **consistent gains with no extra overhead**, which is a **free lunch** for high quality object detection ~
- We hope that **this dynamic viewpoint** can inspire further researches in the future



Source code

Thank you!

Kevin.hkzhang@gmail.com